

# Text Entry Specs

*(by Kaja Mueller-Wang and Duncan Paterson, Cluster Asia and Europe in a Global Context, Heidelberg University, Germany)*

*Version 1.1 April 22, 2015*

## Table of Contents

[File Conventions](#)

[General Markup](#)

[Page Breaks, Page Numbers and Running Heads](#)

[Text Blocks](#)

[Headings](#)

[Paragraphs](#)

[Line Breaks](#)

[Front and Back matter](#)

[Structured Text](#)

[Tables](#)

[Lists](#)

[Table of Contents \(目錄\)](#)

[Special Symbols and Characters](#)

[Character variants and unknown characters](#)

[Special Symbols](#)

[Type Styles](#)

[Small Characters](#)

[Underlinings](#)

[List of All Tags](#)

## File Conventions

Save each book as xml file (.xml) with UTF-8 encoding. Please ensure that all xml files are well-formed. The typed texts should appear inside the root element <text></text>.

Example:

```
<?xml version="1.0" encoding="UTF-8"?>
<text>
  [pdf contents here]
</text>
```

## General Markup

Type the contents of each page separately before moving on to the next page.

### Page Breaks, Page Numbers and Running Heads

Page breaks are marked by “<pb/>”. Running heads appear either in the left or right margin of the page. Type the running heads immediately before or after page breaks. If there is a printed page number on the page include it as part of the running head. Running heads are marked by <fw type="footer"> for running heads appearing in the left margin of the page and <fw type="header"> for running heads appearing in the right margin of the page. If a page number appears within the running head include it between the <fw> tags. Please type large spaces in the running head as a single IDEOGRAPHIC SPACE character U+3000.

Please note: In the digitization of the book, the two half-pages may be on the same scan or on two consecutive scans.

Example:

```
<fw type="footer">文章辨體序說 一三</fw>
```

### Text Blocks

Each section of text that forms a semantic unit on a single or stretching multiple pages should be marked by “<div>” element tags. Headers, titles, subtitles, etc. are part of the <div> and should appear inside them.

### Headings

Enclose each heading (unless it appears in a running head, s.a.) with <head></head> tags. Please note in the books to be typed headings always appear between <div></div> elements (s.b.).

### Paragraphs

The beginning of a new paragraph is either signaled by a header or by indentation. Paragraphs should be marked by <p></p> elements.

## Line Breaks

Type "<lb/>" at the end of each line to determine line-breaks.

Example:

A section beginning with a header followed by paragraphs, including line breaks.

```
<div>
  <head>樂府</head><lb/>
  <p>易曰[ ...]得錄云。</p><lb/>
</div>
```

## Front and Back matter

If prefaces, title pages, table of contents (s.b.), etc, appear in the text they should be enclosed by "<front>" or "<back>" respectively. Front and back matter can appear relative to the whole text, or relative to one of the titles within the volume. If more than one book has been combined within the scans that we submitted, each book forms a "<group>".

Example: 涵芬樓文談

涵芬樓文談 contains two books in one printed volume with prefatory matter for the whole book, as well as prefatory matter of the two books included in this volume. The book corresponds to the <text> elements. <front> covers the preface to the the whole book. The two titles included in the book are marked as a <group>. Immediately followed by <text>. The prefatory matter of the first book appears inside another <front> element. The main body of the text appears inside a <body> element.

In the more common case where only one title appears inside a book, there is a single pair of <body></body> elements to mark the main text, as opposed to dedications, colophons, etc.

```
<text>
  <front>
    <div>
      <head>校點前言</head><lb/>
      <p>[Some text here.]</p>
    </div>
  </front>
  <group>
    <text>
      <front>
        <head>文章辨體序說</head><lb/>
      </front>
      <body>
        <div>
          <head>[Some text here.]</head>
          <p>[Some text here.]</p>
        </div>
      </body>
    </text>
  </text>
```

```
<front>
  <head>文體明辨序說</head><lb/>
</front>
<body>
  <div><p>[Some text here.]</p></div>
</body>
</text>
</group>
</text>
```

## Structured Text

### Tables

Tables are enclosed by `<table></table>` including the number of rows and columns. Include the number of columns as `col="[some number]"` and the number of rows as `row="[some number]"`. The head or title of the table is included as `<head>` (s.a.) inside the table. The sequence of rows and columns is in the order of reading of the main body of text. (From right top to left bottom in the example below). Line-breaks are omitted inside tables.

Example<sup>1</sup>:

A table containing special symbols.

---

<sup>1</sup> For list of the symbols in row two see below [Special Symbols and Characters](#).

|  |  |
|--|--|
| <pre> &lt;table rows="3" cols="9"&gt;   &lt;head&gt;大明唐順之批點法&lt;/head&gt;   &lt;row&gt;     &lt;cell&gt;長國&lt;/cell&gt;     &lt;cell&gt;短間&lt;/cell&gt;     &lt;cell&gt;長點&lt;/cell&gt;     .     .     .   &lt;/row&gt;   &lt;row&gt;     &lt;cell&gt;○○○○○○○○○○&lt;/cell&gt;     &lt;cell&gt;○○&lt;/cell&gt;     &lt;cell&gt;、 、 、 、 、 、 、 &lt;/cell&gt;     .     .     .   &lt;/row&gt;   &lt;row&gt;     &lt;cell&gt;精華&lt;/cell&gt;     &lt;cell&gt;字眼&lt;/cell&gt;     &lt;cell&gt;精華&lt;/cell&gt;     .     .     .   &lt;/row&gt; &lt;/table&gt; </pre> |  |
|--|--|

Lists

Lists should be embedded within “<list>” tags. Please follow the logical order of the list by nesting the individual items to create new lists. Since we follow the logical order within lists, line-breaks are omitted.

Example:  
A generic nested list.

```

<list>
  <item>1</item>
  <item>2
    <list>
      <item>2.1</item>
      <item>2.2</item>
    </list>
  </item>
</list>

```

## Table of Contents (目錄)

Tables of Contents are lists, we are not interested in their layout just their logical structure. Sub-sections should appear within their appropriate sections. Use one HORIZONTAL ELLIPSIS "..." U+2026 for dots that fill spaces in lists

Example:

Sample page of Table of Contents.

|  |   |
|--|---|
| <pre>&lt;div type="contents"&gt;   &lt;head&gt;目錄&lt;/head&gt;   &lt;list&gt;     &lt;item&gt;文章辨體序&lt;hi&gt; (彭時)   &lt;/hi&gt;...七&lt;/item&gt;     &lt;item&gt;文章辨凡例...九&lt;/item&gt;   &lt;/list&gt;   &lt;list&gt;     &lt;item&gt;諸儒總論作文法...一一&lt;/item&gt;   &lt;/list&gt;   &lt;list&gt;     &lt;item&gt;古歌謠詞...一九&lt;/item&gt;     &lt;item&gt;古賦...一九&lt;/item&gt;   &lt;/list&gt;     &lt;item&gt;楚...二〇&lt;/item&gt;     &lt;item&gt;兩漢...二〇&lt;/item&gt;   &lt;/list&gt;     &lt;item&gt;附錄...二一&lt;/item&gt;   &lt;/list&gt;     &lt;item&gt;三國六朝...二一&lt;/item&gt;     &lt;item&gt;唐...二二&lt;/item&gt;     &lt;item&gt;宋...二三&lt;/item&gt;   &lt;/list&gt;   [...]   &lt;fw type="footer"&gt;文章辨體序說 三&lt;/fw&gt;   &lt;pb/&gt;   [...]   &lt;/list&gt;   &lt;/list&gt; &lt;/div&gt;</pre> | <p>目錄</p> <p>文章辨體序 (彭時) ..... 七</p> <p>文章辨體凡例 ..... 九</p> <p>諸儒總論作文法 ..... 一一</p> <p>古歌謠辭 ..... 一九</p> <p>古賦 ..... 一九</p> <p>楚 ..... 二〇</p> <p>兩漢 ..... 二〇</p> <p>附錄 ..... 二一</p> <p>三國六朝 ..... 二一</p> <p>唐 ..... 二二</p> <p>宋 ..... 二三</p> <p>樂府 ..... 二四</p> <p>郊廟歌辭 (吉禮) ..... 二五</p> <p>禮樂歌辭 (軍禮) ..... 二六</p> <p>橫吹曲辭 (原本無, 據「式」增) ..... 二六</p> <p>文章辨體序說 ..... 三</p> |
|--|---|

## Special Symbols and Characters

### Character variants and unknown characters

If a character variant is included in Unicode 7.0, type it. Do not normalize the variant.

Example:

If the character variant 𐀀 of the character 𐀀 occurs in the text, type 𐀀 (U+6B74).

If a character variant is not included in Unicode 7.0, type the standard character instead. In addition, mark it by <g></g>.

Example:

A standard character 𐀀 (U+5602) replacing an obscure non-unicode character “𐀀” in the original text.

<g>𐀀</g>

If there is an unknown character in the text, i.e. a character variant which is readable but where you cannot identify the standard character, add it to the numbered list of unknown characters. For each character that is unknown use a question mark enclosed by <g></g>.

Example:

Three unknown or illegible characters

<g>???</g>

### Special Symbols

| Name                        | Symbol | Unicode | Substitute | Unicode |
|-----------------------------|--------|---------|------------|---------|
| LARGE CIRCLE                | ○      | U+25EF  |            |         |
| HORIZONTAL ELLIPSIS         | ...    | U+2026  |            |         |
| BOX DRAW. DOUBLE HORIZONTAL | =      | U+2550  |            |         |
| BOX DRAW. HEAVY HORIZONTAL  | —      | U+2501  |            |         |
| BOX DRAW. HEAVY VERTICAL    |        | U+2503  |            |         |
| IDEOGRAPHIC SPACE           |        | U+3000  |            |         |
| LESS-THAN                   | <      | U+003C  | <          | U+3008  |
| GREATER-THAN                | >      | U+003E  | >          | U+3009  |

The special characters “<” and “>” are reserved for the markup. Should these appear inside the text, please substitute them with the characters from this list.

## Type Styles

### Small Characters

Strings of small characters, usually used for commentaries, are marked by <hi> </hi>.

Example:

Heading with two nested uses of small script.

|  |                                  |
|--|----------------------------------|
| <pre>&lt;head&gt;古歌謠辭   &lt;hi&gt;（歌、謠、謳、誦、詩、辭、諺附）     &lt;hi&gt;〔注原無，據&lt;line&gt;茅建&lt;/line&gt;本補。〕   &lt;/hi&gt; &lt;/hi&gt; &lt;/head&gt;&lt;lb&gt;</pre> | 古歌謠辭（歌、謠、謳、誦、詩、辭、諺附）〔注原無，據茅建本補。〕 |
|--|----------------------------------|

### Underlinings

A single line next to characters is marked by <line> </line>. A dotted line next to characters is marked by <del> </del>. A wavy line next to characters is marked by <mod> </mod>.

Please note: In old texts, the lines are to the right of characters. In modern texts, the lines may also be to the left of characters. The position to the left or right is not encoded.

As with small characters above, a list of multiple



Example:

Different underlines in text.

|  |                     |  |                                |
|--|---------------------|--|--------------------------------|
| <p>[...]<br/> 後&lt;line&gt;太原&lt;/line&gt;&lt;line&gt;郭茂倩<br/> &lt;/line&gt;輯&lt;mod&gt;樂府&lt;/mod&gt;白卷,<br/> [...]</p> | <p>後太原郭茂倩輯樂府白卷,</p> | <p>[...]<br/> 此無他, &lt;del&gt;無所以然之理以實<br/> 乎其中故也&lt;/del&gt;。大凡一題到<br/> [...]</p> | <p>此無他，無所以然之理以貫乎其中故也。大凡一題到</p> |
|--|---------------------|--|--------------------------------|

Please Note: In the example above there are two single lines and hence two <line></line> elements. Dots on the other hand, only require one pair of <del></del> elements for the whole string of characters.

## List of All Tags

| section   | tag   | name  |
|---|---|---|
| <a href="#">Page Breaks</a>                               | <pb/>   | page break  |
| <a href="#">Running Head</a>                              | <fw type="header"></fw><br><fw type="footer"></fw>                                    | frame work  |
| <a href="#">Text Blocks</a>                               | <div></div>   | division  |
| <a href="#">Headings</a>                                  | <head></head>   | heading   |
| <a href="#">Paragraphs</a>                                | <p></p>   | paragraph   |
| <a href="#">Line Breaks</a>                               | <lb/>   | line break  |
| <a href="#">Front and Back matter</a>                     | <text></text><br><group></group><br><front></front><br><body></body><br><back></back> | text<br>group of texts<br>front matter<br>body<br>back matter |
| <a href="#">Lists</a>                                     | <list><br><item></item><br></list>  | lists with items  |
| <a href="#">Tables</a>                                    | <table row="1" col="1"><br><row><br><cell></cell><br></row><br></table>               | table with 1 row, 1 column, and no head.                      |
| <a href="#">Character variants and unknown characters</a> | <g></g>   | character variants, or unknown characters                     |
| <a href="#">Type Styles</a>                               | <line></line><br><mod></mod><br><del></del>   | single line<br>wavy line<br>dotted line                       |
| <a href="#">Small Characters</a>                          | <hi></hi>   | small type / commentary                                       |